

Error Correction for Index Coding with Side Information

Son Hoang Dau^{*}, Vitaly Skachek^{†,1}, and Yeow Meng Chee[‡]

^{*,‡}Division of Mathematical Sciences, School of Physical and Mathematical Sciences
Nanyang Technological University, 21 Nanyang Link, Singapore 637371

[†]Coordinated Science Laboratory, University of Illinois at Urbana-Champaign
1308 W. Main Street, Urbana, IL 61801, USA

Emails: ^{*}daus0002@ntu.edu.sg, [†]vitalys@illinois.edu, [‡]YMChae@ntu.edu.sg

Abstract—A problem of index coding with side information was first considered by Y. Birk and T. Kol (*IEEE INFOCOM*, 1998). In the present work, a generalization of index coding scheme, where transmitted symbols are subject to errors, is studied. Error-correcting methods for such a scheme, and their parameters, are investigated. In particular, the following question is discussed: given the side information hypergraph of index coding scheme and the maximal number of erroneous symbols δ , what is the shortest length of a linear index code, such that every receiver is able to recover the required information? This question turns out to be a generalization of the problem of finding a shortest-length error-correcting code with a prescribed error-correcting capability in the classical coding theory.

The Singleton bound and two other bounds, referred to as the α -bound and the κ -bound, for the optimal length of a linear error-correcting index code (ECIC) are established. For large alphabets, a construction based on concatenation of an optimal index code with an MDS classical code, is shown to attain the Singleton bound. For smaller alphabets, however, this construction may not be optimal. A random construction is also analyzed. It yields another inexplicit bound on the length of an optimal linear ECIC.

Further, the problem of error-correcting decoding by a linear ECIC is studied. It is shown that in order to decode correctly the desired symbol, the decoder is required to find one of the vectors, belonging to an affine space containing the actual error vector. The syndrome decoding is shown to produce the correct output if the weight of the error pattern is less or equal to the error-correcting capability of the corresponding ECIC.

Finally, the notion of static ECIC, which is suitable for use with a family of instances of an index coding problem, is introduced. Several bounds on the length of static ECIC's are derived, and constructions for static ECIC's are discussed. Connections of these codes to weakly resilient Boolean functions are established.

Index Terms—index coding, network coding, side information, error correction, minimum distance, broadcast.

I. INTRODUCTION

A. Background

¹The work of this author was done while he was with the Division of Mathematical Sciences, School of Physical and Mathematical Sciences, Nanyang Technological University, 21 Nanyang Link, Singapore 637371.

A part of this work is to be presented in the *IEEE International Symposium on Information Theory (ISIT)*, St. Petersburg, Russia, July-August 2011.

The problem of Index Coding with Side Information (ICSI) was introduced by Birk and Kol [1], [2]. During the transmission, each client might miss a certain part of the data, due to intermittent reception, limited storage capacity or any other reasons. Via a slow backward channel, the clients let the server know which messages they already have in their possession, and which messages they are interested to receive. The server has to find a way to deliver to each client all the messages he requested, yet spending a minimum number of transmissions. As it was shown in [1], the server can significantly reduce the number of transmissions by coding the messages.

The toy example in Figure 1 presents a scenario with one broadcast transmitter and four receivers. Each receiver requires a different information packet (we sometimes simply call it message). The naïve approach requires four separate transmissions, one transmission per an information packet. However, by exploiting the knowledge on the subsets of messages that clients already have, and by using coding of the transmitted data, the server can just broadcast one coded packet.

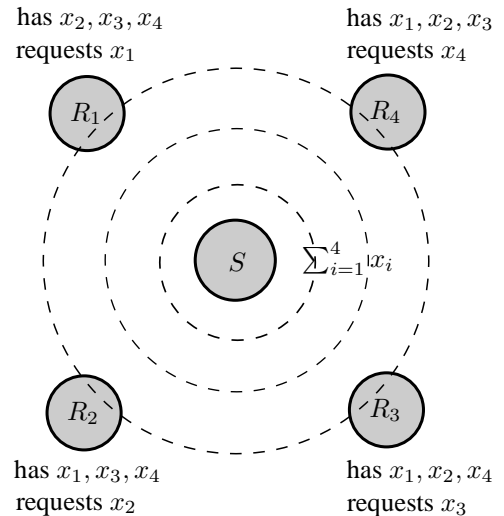


Fig. 1: An example of the ICSI problem

Possible applications of index coding include communica-

tions scenarios, in which a satellite or a server broadcasts a set of messages to a set clients, such as daily newspaper delivery or video-on-demand. Index coding with side information can also be used in opportunistic wireless networks. These are the networks in which a wireless node can opportunistically listen to the wireless channel. The client may obtain packets that are not designated to it (see [3]–[5]). As a result, a node obtains some side information about the transmitted data. Exploiting this additional knowledge may help to increase the throughput of the system.

The ICSI problem has been a subject of several recent studies [3], [6]–[13]. This problem can be viewed as a special case of the Network Coding (NC) problem [14], [15]. In particular, as it was shown in [3], [11], every instance of the NC problem can be reduced to an instance of the ICSI problem.

B. Our contribution

The preceding works on the ICSI problem consider scenario where the transmissions are error-free. In practice, of course, this might not be the case. In this work, we assume that the transmitted symbols are subject to errors. We extend some known results on index coding to a case where any receiver can correct up to a certain number of errors. It turns out that the problem of designing such error-correcting index codes (ECIC's) naturally generalizes the problem of constructing classical error-correcting codes.

More specifically, assume that the number of messages that the server possesses is n , and that the designed maximal number of errors is δ . We show that the problem of constructing ECIC of minimal possible length is equivalent to the problem of constructing a matrix \mathbf{L} which has n rows and the minimal possible number of columns, such that

$$\text{wt}(\mathbf{z}\mathbf{L}) \geq 2\delta + 1 \text{ for all } \mathbf{z} \in \mathcal{I},$$

where \mathcal{I} is a certain subset of $\mathbb{F}_q^n \setminus \{\mathbf{0}\}$. Here $\text{wt}(\mathbf{x})$ denotes the Hamming weight of the vector \mathbf{x} , \mathbb{F}_q stands for a finite field with q elements, and $\mathbf{0}$ is the all-zeros vector. If $\mathcal{I} = \mathbb{F}_q^n \setminus \{\mathbf{0}\}$, this problem becomes equivalent to the problem of designing a shortest-length linear code of given dimension and minimum distance.

In this work, we establish an upper bound (the κ -bound) and a lower bound (the α -bound) on the shortest length of a linear ECIC, which is able to correct any error pattern of size up to δ . More specifically, let \mathcal{H} be the side information hypergraph that describes the instance of the ICSI problem. Let $\mathcal{N}_q[\mathcal{H}, \delta]$ denote the length of a shortest-length linear ECIC over \mathbb{F}_q , such that every R_i can recover the desired message, if the number of errors is at most δ . We use notation $N_q[k, d]$ for the length of an optimal linear error-correcting code of dimension k and minimum distance d over \mathbb{F}_q . We obtain

$$N_q[\alpha(\mathcal{H}), 2\delta + 1] \leq \mathcal{N}_q[\mathcal{H}, \delta] \leq N_q[\kappa_q(\mathcal{H}), 2\delta + 1], \quad (1)$$

where $\alpha(\mathcal{H})$ is the generalized independence number and $\kappa_q(\mathcal{H})$ is the min-rank (over \mathbb{F}_q) of \mathcal{H} .

For linear index codes, we also derive an analog of the Singleton bound. This result implies that (over sufficiently

large alphabet) the concatenation of a standard MDS error-correcting code with an optimal linear index code yields an optimal linear error-correcting index code. Finally, we consider random ECIC's. By analyzing its parameters, we obtain an upper bound on its length.

When the side information hypergraph is a pentagon, and $\delta = 2$, the inequalities in (1) are shown to be strict. This implies that a concatenated scheme based on a classical error-correcting code and on a linear non-error-correcting index code does not necessarily yield an optimal linear error-correcting index code. Since ICSI problem can also be viewed as a source coding problem [6], [13], this example demonstrates that sometimes designing a single code for both source and channel coding can result in a smaller number of transmissions.

The decoding of a linear ECIC is somewhat different from that of a classical error-correcting code. There is no longer a need for a complete recovery of the whole information vector. We analyze the decoding criteria for the ECIC's and show that the syndrome decoding, which might be different for each receiver, results in a correct result, provided that the number of errors does not exceed the error-correcting capability of the code.

An ECIC is called static under a family of instances of the ICSI problem if it works for all of these instances. Such an ECIC is interesting since it remains useful as long as the parameters of the problem vary within a particular range. Bounds and constructions for static ECIC's are studied in Section VIII. Connections between static ECIC's and weakly resilient vectorial Boolean functions are also discussed.

The problem of error correction for NC was studied in several previous works. However, these results are not directly applicable to the ICSI problem. First, there is only a very limited variety of results for non-multicast networks in the existing literature. The ICSI problem, however, is a special case of the non-multicast NC problem. Second, the ICSI problem can be modeled by the NC scenario [3], yet, this requires that there are directed edges from particular sources to each sink, which provide the side information. The symbols transmitted on these special edges are not allowed to be corrupted. By contrast, for error-correcting NC, symbols transmitted on all edges can be corrupted.

The paper is organized as follows. Basic notations and definitions, used throughout the paper, are provided in Section II. The problem of index coding with and without error-correction is introduced in Section III. Some basic results are presented in that section. The α -bound and the κ -bound are derived in Section IV. The Singleton bound is presented in Section V. Random codes are discussed in Section VI. Syndrome decoding is studied in Section VII. A notion of static error-correcting index codes is presented in Section VIII. Several bounds on the length of such codes are derived, and connections to resilient function are shown in that section. Finally, the results are summarized in Section IX, and some open questions are proposed therein.

II. PRELIMINARIES

In this section we introduce some useful notation. Here \mathbb{F}_q is the finite field of q elements, where q is a power of prime,

and \mathbb{F}_q^* is the set of all nonzero elements of \mathbb{F}_q .

Let $[n] = \{1, 2, \dots, n\}$. For the vectors $\mathbf{u} = (u_1, u_2, \dots, u_n) \in \mathbb{F}_q^n$ and $\mathbf{v} = (v_1, v_2, \dots, v_n) \in \mathbb{F}_q^n$, the (Hamming) distance between \mathbf{u} and \mathbf{v} is defined to be the number of coordinates where \mathbf{u} and \mathbf{v} differ, namely,

$$d(\mathbf{u}, \mathbf{v}) = |\{i \in [n] : u_i \neq v_i\}|.$$

If $\mathbf{u} \in \mathbb{F}_q^n$ and $M \subseteq \mathbb{F}_q^n$ is a set of vectors (or a vector subspace), then the last definition can be extended to

$$d(\mathbf{u}, M) = \min_{\mathbf{v} \in M} d(\mathbf{u}, \mathbf{v}).$$

The *support* of a vector $\mathbf{u} \in \mathbb{F}_q^n$ is defined to be the set $\text{supp}(\mathbf{u}) = \{i \in [n] : u_i \neq 0\}$. The (Hamming) weight of a vector \mathbf{u} , denoted $\text{wt}(\mathbf{u})$, is defined to be $|\text{supp}(\mathbf{u})|$, the number of nonzero coordinates of \mathbf{u} . Suppose $E \subseteq [n]$. We write $\mathbf{u} \triangleleft E$ whenever $\text{supp}(\mathbf{u}) \subseteq E$.

A k -dimensional subspace \mathcal{C} of \mathbb{F}_q^n is called a linear $[n, k, d]_q$ code over \mathbb{F}_q if the minimum distance of \mathcal{C} ,

$$d(\mathcal{C}) \triangleq \min_{\mathbf{u} \in \mathcal{C}, \mathbf{v} \in \mathcal{C}, \mathbf{u} \neq \mathbf{v}} d(\mathbf{u}, \mathbf{v}),$$

is equal to d . Sometimes we may use the notation $[n, k]_q$ for the sake of simplicity. The vectors in \mathcal{C} are called codewords. It is easy to see that the minimum weight of a nonzero codeword in a linear code \mathcal{C} is equal to its minimum distance $d(\mathcal{C})$. A *generator matrix* \mathbf{G} of an $[n, k]_q$ code \mathcal{C} is a $k \times n$ matrix whose rows are linearly independent codewords of \mathcal{C} . Then $\mathcal{C} = \{\mathbf{y}\mathbf{G} : \mathbf{y} \in \mathbb{F}_q^k\}$. The *parity-check matrix* of \mathcal{C} is an $(n-k) \times n$ matrix \mathbf{H} over \mathbb{F}_q such that $\mathbf{c} \in \mathcal{C} \Leftrightarrow \mathbf{H}\mathbf{c}^T = \mathbf{0}^T$. Given q, k , and d , let $N_q[k, d]$ denote the length of the shortest linear code over \mathbb{F}_q which has dimension k and minimum distance d .

We use $\mathbf{e}_i = (\underbrace{0, \dots, 0}_{i-1}, 1, \underbrace{0, \dots, 0}_{n-i}) \in \mathbb{F}_q^n$ to denote the unit vector, which has a one at the i th position, and zeros elsewhere. For a vector $\mathbf{y} = (y_1, y_2, \dots, y_n)$ and a subset $B = \{i_1, i_2, \dots, i_b\}$ of $[n]$, where $i_1 < i_2 < \dots < i_b$, let \mathbf{y}_B denote the vector $(y_{i_1}, y_{i_2}, \dots, y_{i_b})$.

For an $n \times N$ matrix \mathbf{L} , let \mathbf{L}_i denote its i th row. For a set $E \subseteq [n]$, let \mathbf{L}_E denote the $|E| \times N$ matrix obtained from \mathbf{L} by deleting all the rows of \mathbf{L} which are not indexed by the elements of E . For a set of vectors M , we use notation $\text{span}(M)$ to denote the linear space spanned by the vectors in M . We also use notation $\text{colspan}(\mathbf{L})$ for the linear space spanned by the columns of the matrix \mathbf{L} .

Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ be a graph with a vertex set \mathcal{V} and an edge set \mathcal{E} . The graph is called *undirected* if every edge $e \in \mathcal{E}$, $e = \{u, v\}$, and $u, v \in \mathcal{V}$. A graph \mathcal{G} is *directed* if every edge $e \in \mathcal{E}$ is an ordered pair $e = (u, v)$, $u, v \in \mathcal{V}$. A directed graph \mathcal{G} is called *symmetric* if

$$(u, v) \in \mathcal{E} \Leftrightarrow (v, u) \in \mathcal{E}.$$

There is a natural correspondence between undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ and directed symmetric graph $\mathcal{G}' = (\mathcal{V}, \mathcal{E}')$ defined as

$$\mathcal{E} = \{\{u, v\} : (u, v) \in \mathcal{E}'\}. \quad (2)$$

Let \mathcal{G} be an undirected graph. A subset of vertices $\mathcal{S} \subseteq \mathcal{V}$ is called an *independent set* if $\forall u, v \in \mathcal{S}$, $\{u, v\} \notin \mathcal{E}$. The size

of the largest independent set in \mathcal{G} is called the *independence number* of \mathcal{G} , and is denoted by $\alpha(\mathcal{G})$. The graph $\bar{\mathcal{G}} = (\mathcal{V}, \bar{\mathcal{E}})$ is called the *complement* of $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ if

$$\bar{\mathcal{E}} = \{\{u, v\} : u \in \mathcal{V}, v \in \mathcal{V}, \{u, v\} \notin \mathcal{E}\}.$$

A *coloring* of \mathcal{G} using χ colors is a function $\psi : \mathcal{V} \rightarrow [\chi]$, such that

$$\forall e = \{u, v\} \in \mathcal{E} : \psi(u) \neq \psi(v).$$

The *chromatic number* of \mathcal{G} is the smallest number χ such that there exists a coloring of \mathcal{G} using χ colors, and it is denoted by $\chi(\mathcal{G})$. By using the correspondence (2), the definitions of independence number, graph complement and chromatic number are trivially extended to directed symmetric graphs.

III. INDEX CODING AND ERROR CORRECTION

A. Index Coding with Side Information

Index Coding with Side Information problem considers the following communications scenario. There is a unique sender (or source) S , who has a vector of messages $\mathbf{x} = (x_1, x_2, \dots, x_n)$ in his possession. There are also m receivers R_1, R_2, \dots, R_m , receiving information from S via a broadcast channel. For each $i \in [m]$, R_i has side information, i.e. R_i owns a subset of messages $\{x_j\}_{j \in \mathcal{X}_i}$, where $\mathcal{X}_i \subseteq [n]$. Each R_i , $i \in [m]$, is interested in receiving the message $x_{f(i)}$ (we say that R_i requires $x_{f(i)}$), where the mapping $f : [m] \rightarrow [n]$ satisfies $f(i) \notin \mathcal{X}_i$ for all $i \in [m]$. Hereafter, we use the notation $\mathcal{X} = (\mathcal{X}_1, \mathcal{X}_2, \dots, \mathcal{X}_m)$. An instance of the ICSI problem is given by a quadruple (m, n, \mathcal{X}, f) . It can also be conveniently described by a directed hypergraph [13].

Definition 3.1: Let (m, n, \mathcal{X}, f) be an instance of the ICSI problem. The corresponding *side information (directed) hypergraph* $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$ is defined by the vertex set $\mathcal{V} = [n]$ and the edge set $\mathcal{E}_{\mathcal{H}}$, where

$$\mathcal{E}_{\mathcal{H}} = \{(f(i), \mathcal{X}_i) : i \in [m]\}.$$

We often refer to (m, n, \mathcal{X}, f) as an instance of the ICSI problem described by the hypergraph \mathcal{H} .

Each side information hypergraph $\mathcal{H} = (\mathcal{V}, \mathcal{E}_{\mathcal{H}})$ can be associated with the directed graph $\mathcal{G}_{\mathcal{H}} = (\mathcal{V}, \mathcal{E})$ in the following way. For each directed edge $(f(i), \mathcal{X}_i) \in \mathcal{E}_{\mathcal{H}}$ there will be $|\mathcal{X}_i|$ directed edges $(f(i), v) \in \mathcal{E}$, for $v \in \mathcal{X}_i$. When $m = n$ and $f(i) = i$ for all $i \in [m]$, the graph $\mathcal{G}_{\mathcal{H}}$ is, in fact, the *side information graph*, defined in [6].

The goal of the ICSI problem is to design a coding scheme that allows S to satisfy the requests of all receivers R_i in the least number of transmissions. More formally, we have the following definition.

Definition 3.2: An *index code* over \mathbb{F}_q for an instance of the ICSI problem described by $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$ (or just an \mathcal{H} -IC over \mathbb{F}_q), is an encoding function

$$\mathfrak{E} : \mathbb{F}_q^n \rightarrow \mathbb{F}_q^N,$$

such that for each receiver R_i , $i \in [m]$, there exists a decoding function

$$\mathfrak{D}_i : \mathbb{F}_q^N \times \mathbb{F}_q^{|\mathcal{X}_i|} \rightarrow \mathbb{F}_q ,$$

satisfying

$$\forall \mathbf{x} \in \mathbb{F}_q^n : \mathfrak{D}_i(\mathfrak{E}(\mathbf{x}), \mathbf{x}_{\mathcal{X}_i}) = x_{f(i)} .$$

Sometimes we refer to such \mathfrak{E} as a *non-error-correcting* index code. The parameter N is called the *length* of the index code. In the scheme corresponding to this code, S broadcasts a vector $\mathfrak{E}(\mathbf{x})$ of length N over \mathbb{F}_q .

Definition 3.3: A *linear index code* is an index code, for which the encoding function \mathfrak{E} is a linear transformation over \mathbb{F}_q . Such a code can be described as

$$\forall \mathbf{x} \in \mathbb{F}_q^n : \mathfrak{E}(\mathbf{x}) = \mathbf{x}\mathbf{L} ,$$

where \mathbf{L} is an $n \times N$ matrix over \mathbb{F}_q . The matrix \mathbf{L} is called the *matrix corresponding to the index code* \mathfrak{E} . The code \mathfrak{E} is also referred to as the *linear index code based on \mathbf{L}* .

Hereafter, we assume that $\mathcal{X} = (\mathcal{X}_i)_{i \in [m]}$ is known to S . Moreover, we also assume that the code \mathfrak{E} is known to each receiver R_i , $i \in [m]$. In practice this can be achieved by a preliminary communication session, when the knowledge of the sets \mathcal{X}_i for $i \in [m]$ and of the code \mathfrak{E} are disseminated between the participants of the scheme.

Definition 3.4: Suppose $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$ corresponds to an instance of the ICSI problem. Then the *min-rank* of \mathcal{H} over \mathbb{F}_q is defined as

$$\kappa_q(\mathcal{H}) \triangleq \min\{\text{rank}_{\mathbb{F}_q}(\{\mathbf{v}_i + \mathbf{e}_{f(i)}\}_{i \in [m]}) : \mathbf{v}_i \in \mathbb{F}_q^n, \mathbf{v}_i \triangleleft \mathcal{X}_i\} .$$

Observe that $\kappa_q(\mathcal{H})$ generalizes the min-rank over \mathbb{F}_q of the side information graph, which was defined in [6]. More specifically, when $m = n$ and $f(i) = i$ for all $i \in [m]$, $\mathcal{G}_{\mathcal{H}}$ becomes the side information graph, and $\kappa_q(\mathcal{H}) = \text{min-rank}_q(\mathcal{G}_{\mathcal{H}})$. The min-rank of an undirected graph was first introduced by Haemers [16] to bound the Shannon capacity of a graph, and was later proved in [6], [7] to be the smallest number of transmissions in a linear index code.

The following lemma was implicitly formulated in [6] for the case where $m = n$, $q = 2$, $f(i) = i$ for all $i \in [n]$, and generalized to its current form in [17].

Lemma 3.5: Consider an instance of the ICSI problem described by $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$.

- 1) The matrix \mathbf{L} corresponds to a linear \mathcal{H} -IC over \mathbb{F}_q if and only if for each $i \in [m]$ there exists $\mathbf{v}_i \in \mathbb{F}_q^n$ such that
 - $\mathbf{v}_i \triangleleft \mathcal{X}_i$;
 - $\mathbf{v}_i + \mathbf{e}_{f(i)} \in \text{colspan}(\mathbf{L})$.
- 2) The smallest possible length of a linear \mathcal{H} -IC over \mathbb{F}_q is $\kappa_q(\mathcal{H})$.

B. Error-Correcting Index Code with Side Information

Due to noise, the symbols received by R_i , $i \in [m]$, may be subject to errors. Consider an ICSI instance (m, n, \mathcal{X}, f) , and assume that S broadcasts a vector $\mathfrak{E}(\mathbf{x}) \in \mathbb{F}_q^N$. Let $\epsilon_i \in \mathbb{F}_q^N$ be the error affecting the information received by R_i , $i \in [m]$. Then R_i actually receives the vector

$$\mathbf{y}_i = \mathfrak{E}(\mathbf{x}) + \epsilon_i \in \mathbb{F}_q^N ,$$

instead of $\mathfrak{E}(\mathbf{x})$. The following definition is a generalization of Definition 3.2.

Definition 3.6: Consider an instance of the ICSI problem described by $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$. A δ -*error-correcting index code* $((\delta, \mathcal{H})\text{-ECIC})$ over \mathbb{F}_q for this instance is an encoding function

$$\mathfrak{E} : \mathbb{F}_q^n \rightarrow \mathbb{F}_q^N ,$$

such that for each receiver R_i , $i \in [m]$, there exists a decoding function

$$\mathfrak{D}_i : \mathbb{F}_q^N \times \mathbb{F}_q^{|\mathcal{X}_i|} \rightarrow \mathbb{F}_q ,$$

satisfying

$$\forall \mathbf{x}, \epsilon_i \in \mathbb{F}_q^n, \text{wt}(\epsilon_i) \leq \delta : \mathfrak{D}_i(\mathfrak{E}(\mathbf{x}) + \epsilon_i, \mathbf{x}_{\mathcal{X}_i}) = x_{f(i)} .$$

The definitions of the length, of a linear index code, and of the matrix corresponding to an index code are naturally extended to an error-correcting index code. Note that if \mathfrak{E} is an \mathcal{H} -IC, then it is a $(0, \mathcal{H})\text{-ECIC}$, and vice versa.

Definition 3.7: An *optimal* linear $(\delta, \mathcal{H})\text{-ECIC}$ over \mathbb{F}_q is a linear $(\delta, \mathcal{H})\text{-ECIC}$ over \mathbb{F}_q of the smallest possible length $\mathcal{N}_q[\mathcal{H}, \delta]$.

Consider an instance of the ICSI problem described by $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$. We define the set of vectors

$$\mathcal{I}(q, \mathcal{H}) \triangleq \{\mathbf{z} \in \mathbb{F}_q^n : \exists i \in [m] \text{ such that } \mathbf{z}_{\mathcal{X}_i} = \mathbf{0} \text{ and } z_{f(i)} \neq 0\} .$$

For all $i \in [m]$, we also define

$$\mathcal{Y}_i \triangleq [n] \setminus (\{f(i)\} \cup \mathcal{X}_i) .$$

Then the collection of supports of all vectors in $\mathcal{I}(q, \mathcal{H})$ is given by

$$\mathcal{J}(\mathcal{H}) \triangleq \bigcup_{i \in [m]} \{\{f(i)\} \cup Y_i : Y_i \subseteq \mathcal{Y}_i\} . \quad (3)$$

The necessary and sufficient condition for a matrix \mathbf{L} to be the matrix corresponding to some $(\delta, \mathcal{H})\text{-ECIC}$ is given in the following lemma.

Lemma 3.8: The matrix \mathbf{L} corresponds to a $(\delta, \mathcal{H})\text{-ECIC}$ over \mathbb{F}_q if and only if

$$\text{wt}(\mathbf{z}\mathbf{L}) \geq 2\delta + 1 \text{ for all } \mathbf{z} \in \mathcal{I}(q, \mathcal{H}) . \quad (4)$$

Equivalently, \mathbf{L} corresponds to a $(\delta, \mathcal{H})\text{-ECIC}$ over \mathbb{F}_q if and only if

$$\text{wt}\left(\sum_{i \in K} z_i \mathbf{L}_i\right) \geq 2\delta + 1 ,$$

for all $K \in \mathcal{J}(\mathcal{H})$ and for all choices of $z_i \in \mathbb{F}_q^*$, $i \in K$.

Proof: For each $\mathbf{x} \in \mathbb{F}_q^n$, we define

$$B(\mathbf{x}, \delta) = \{\mathbf{y} \in \mathbb{F}_q^N : \mathbf{y} = \mathbf{x}\mathbf{L} + \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \in \mathbb{F}_q^N, \text{wt}(\boldsymbol{\epsilon}) \leq \delta\},$$

the set of all vectors resulting from at most δ errors in the transmitted vector associated with the information vector \mathbf{x} . Then the receiver R_i can recover $x_{f(i)}$ correctly if and only if

$$B(\mathbf{x}, \delta) \cap B(\mathbf{x}', \delta) = \emptyset,$$

for every pair $\mathbf{x}, \mathbf{x}' \in \mathbb{F}_q^n$ satisfying:

$$\mathbf{x}_{\mathcal{X}_i} = \mathbf{x}'_{\mathcal{X}_i} \text{ and } x_{f(i)} \neq x'_{f(i)}.$$

(Observe that R_i is interested only in the bit $x_{f(i)}$, not in the whole vector \mathbf{x} .)

Therefore, \mathbf{L} corresponds to a (δ, \mathcal{H}) -ECIC if and only if the following condition is satisfied: for all $i \in [m]$ and for all $\mathbf{x}, \mathbf{x}' \in \mathbb{F}_q^n$ such that $\mathbf{x}_{\mathcal{X}_i} = \mathbf{x}'_{\mathcal{X}_i}$ and $x_{f(i)} \neq x'_{f(i)}$, it holds

$$\forall \boldsymbol{\epsilon}, \boldsymbol{\epsilon}' \in \mathbb{F}_q^N, \text{wt}(\boldsymbol{\epsilon}) \leq \delta, \text{wt}(\boldsymbol{\epsilon}') \leq \delta : \mathbf{x}\mathbf{L} + \boldsymbol{\epsilon} \neq \mathbf{x}'\mathbf{L} + \boldsymbol{\epsilon}'. \quad (5)$$

Denote $\mathbf{z} = \mathbf{x}' - \mathbf{x}$. Then, the condition in (5) can be reformulated as follows: for all $i \in [m]$ and for all $\mathbf{z} \in \mathbb{F}_q^n$ such that $\mathbf{z}_{\mathcal{X}_i} = \mathbf{0}$ and $z_{f(i)} \neq 0$, it holds

$$\forall \boldsymbol{\epsilon}, \boldsymbol{\epsilon}' \in \mathbb{F}_q^N, \text{wt}(\boldsymbol{\epsilon}) \leq \delta, \text{wt}(\boldsymbol{\epsilon}') \leq \delta : \mathbf{z}\mathbf{L} \neq \boldsymbol{\epsilon} - \boldsymbol{\epsilon}'. \quad (6)$$

The equivalent condition is that for all $\mathbf{z} \in \mathcal{I}(q, \mathcal{H})$,

$$\text{wt}(\mathbf{z}\mathbf{L}) \geq 2\delta + 1.$$

Since for $\mathbf{z} \in \mathcal{I}(q, \mathcal{H})$ we have

$$\mathbf{z}\mathbf{L} = \sum_{i \in \text{supp}(\mathbf{z})} z_i \mathbf{L}_i,$$

the condition (4) can be restated as

$$\text{wt}\left(\sum_{i \in K} z_i \mathbf{L}_i\right) \geq 2\delta + 1,$$

for all $K \in \mathcal{J}(\mathcal{H})$ and for all choices of nonzero $z_i \in \mathbb{F}_q$, $i \in K$. ■

The next corollary follows from Lemma 3.8 in a straightforward manner. It is not hard to see that the conditions stated in Lemma 3.8 and in the corollary below are, in fact, equivalent.

Corollary 3.9: For all $i \in [m]$, let

$$\mathbf{M}_i \triangleq \text{span}(\{\mathbf{L}_j : j \in \mathcal{Y}_i\}).$$

Then, the matrix \mathbf{L} corresponds to a (δ, \mathcal{H}) -ECIC over \mathbb{F}_q if and only if

$$\forall i \in [m] : d(\mathbf{L}_{f(i)}, \mathbf{M}_i) \geq 2\delta + 1. \quad (7)$$

The next corollary also follows directly from Lemma 3.8 by considering an error-free setup, i.e. $\delta = 0$. It is easy to verify that the conditions stated in this corollary and in Lemma 3.5 are equivalent, as expected.

Corollary 3.10: The matrix \mathbf{L} corresponds to an \mathcal{H} -IC over \mathbb{F}_q if and only if

$$\text{wt}\left(\sum_{i \in K} z_i \mathbf{L}_i\right) \geq 1,$$

for all $K \in \mathcal{J}(\mathcal{H})$ and for all choices of $z_i \in \mathbb{F}_q^*$, $i \in K$, or, equivalently,

$$\forall i \in [m] : \mathbf{L}_{f(i)} \notin \text{span}(\{\mathbf{L}_j\}_{j \in \mathcal{Y}_i}).$$

Example 3.11: Let $q = 2$, $m = n = 3$, and $f(i) = i$ for $i \in [3]$. Suppose $\mathcal{X}_1 = \{2, 3\}$, $\mathcal{X}_2 = \{1, 3\}$, and $\mathcal{X}_3 = \{1, 2\}$. Let

$$\mathbf{L} = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 1 & 1 & 0 & 1 \\ 1 & 0 & 1 & 1 \end{pmatrix}.$$

Note that \mathbf{L} generates a $[4, 3, 1]_2$ code, which has minimum distance one. However, the index code based on \mathbf{L} can still correct one error. Indeed, let $\mathcal{H} = \mathcal{H}(3, 3, \mathcal{X}, f)$, we have

$$\mathcal{I}(2, \mathcal{H}) = \{100, 010, 001\}.$$

Since each row of \mathbf{L} has weight at least three, it follows that $\text{wt}(\mathbf{z}\mathbf{L}) \geq 3$ for all $\mathbf{z} \in \mathcal{I}(2, \mathcal{H})$. By Lemma 3.8, \mathbf{L} corresponds to a $(1, \mathcal{H})$ -ECIC over \mathbb{F}_2 .

In fact, for this instance, even a simpler index code of length three, based on

$$\mathbf{L}' = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix},$$

is a $(1, \mathcal{H})$ -ECIC over \mathbb{F}_2 .

Example 3.12: Assume that $m = n$ and $f(i) = i$ for all $i \in [m]$. Furthermore, suppose that $\mathcal{X}_i = \emptyset$ for all $i \in [m]$ (i.e. there is no side information available to the receivers). Let $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$. Then, $\mathcal{I}(q, \mathcal{H}) = \mathbb{F}_q^n \setminus \{\mathbf{0}\}$. Hence, by Lemma 3.8, the $n \times N$ matrix \mathbf{L} corresponding to a (δ, \mathcal{H}) -ECIC over \mathbb{F}_q (for some integer $\delta \geq 0$) is a generating matrix of an $[N, n, \geq 2\delta + 1]_q$ linear code. Thus, under these settings, the problem of designing an optimal ECIC is reduced to the problem of constructing an optimal classical linear error-correcting code.

Observe however, that for general \mathcal{X} , changing the order of rows in \mathbf{L} can lead to ECIC's with different error-correcting capabilities. Therefore, the problem of designing an optimal linear ECIC is essentially the problem of finding the matrix \mathbf{L} corresponding to that code. However, the minimum distance of the code generated by the rows of \mathbf{L} is not necessary a valid indicator for goodness of an ECIC. Sometimes, as Example 3.11 shows, matrix \mathbf{L} with redundant rows yields a good ECIC.

IV. THE α -BOUND AND THE κ -BOUND

Let (m, n, \mathcal{X}, f) be an instance of the ICSI problem, and let \mathcal{H} be the corresponding side information hypergraph. Next, we introduce the following definitions for the hypergraph \mathcal{H} .

Definition 4.1: A subset H of $[n]$ is called a *generalized independent set* in \mathcal{H} if every nonempty subset K of H belongs to $\mathcal{J}(\mathcal{H})$.

Definition 4.2: A generalized independent set of the largest size in \mathcal{H} is called a *maximum generalized independent set*. The size of a maximum generalized independent set in \mathcal{H} is called the *generalized independence number*, and denoted by $\alpha(\mathcal{H})$.

When $m = n$ and $f(i) = i$ for all $i \in [n]$, the generalized independence number of \mathcal{H} is equal to the maximum size of an acyclic induced subgraph of $\mathcal{G}_{\mathcal{H}}$, which was introduced in [6]. In particular, when $\mathcal{G}_{\mathcal{H}}$ is symmetric, $\alpha(\mathcal{H})$ is the independence number of $\mathcal{G}_{\mathcal{H}}$. We prove the latter statement in the Appendix.

Next, we present a lower bound on the length of a (δ, \mathcal{H}) -ECIC. We call this bound α -bound.

Theorem 4.3 (α -bound): The length of an optimal linear (δ, \mathcal{H}) -ECIC over \mathbb{F}_q satisfies

$$N_q[\mathcal{H}, \delta] \geq N_q[\alpha(\mathcal{H}), 2\delta + 1].$$

Moreover, the equality is attained if there exists an $n \times \alpha(\mathcal{H})$ matrix $\mathbf{B} = (b_{i,j})$ over \mathbb{F}_q satisfying the following condition: for all $K \in \mathcal{J}(\mathcal{H})$ and for all choices of $z_i \in \mathbb{F}_q^*$, $i \in K$, there always exists some j such that

$$\sum_{i \in K} z_i b_{i,j} \neq 0.$$

Proof: Consider an $n \times N$ matrix \mathbf{L} , which corresponds to a (δ, \mathcal{H}) -ECIC. Let $H = \{i_1, i_2, \dots, i_{\alpha(\mathcal{H})}\}$ be a maximum generalized independent set in \mathcal{H} . Then, every subset $K \subseteq H$ satisfies $K \in \mathcal{J}(\mathcal{H})$. Therefore,

$$\text{wt} \left(\sum_{i \in K} z_i \mathbf{L}_i \right) \geq 2\delta + 1$$

for all $K \subseteq H$, $K \neq \emptyset$, and for all choices of $z_i \in \mathbb{F}_q^*$, $i \in K$. Hence, the $\alpha(\mathcal{H})$ rows of \mathbf{L} , namely $\mathbf{L}_{i_1}, \mathbf{L}_{i_2}, \dots, \mathbf{L}_{i_{\alpha(\mathcal{H})}}$, form a generator matrix of an $[N, \alpha(\mathcal{H}), 2\delta + 1]_q$ code. Therefore,

$$N \geq N_q[\alpha(\mathcal{H}), 2\delta + 1].$$

Next, we assume the existence of a matrix \mathbf{B} satisfying the properties stated in the theorem. Let \mathbf{L}' be a generator matrix of some $[N', \alpha(\mathcal{H}), 2\delta + 1]_q$ code, where $N' = N_q[\alpha(\mathcal{H}), 2\delta + 1]$. We construct the $n \times N'$ matrix \mathbf{L} as follows. For $i \in [n]$, let

$$\mathbf{L}_i = \sum_{j=1}^{\alpha(\mathcal{H})} b_{i,j} \mathbf{L}'_j.$$

For every $K \in \mathcal{J}(\mathcal{H})$ and for all choices of $z_i \in \mathbb{F}_q^*$, $i \in K$, we have

$$\begin{aligned} \text{wt} \left(\sum_{i \in K} z_i \mathbf{L}_i \right) &= \text{wt} \left(\sum_{i \in K} z_i \sum_{j=1}^{\alpha(\mathcal{H})} b_{i,j} \mathbf{L}'_j \right) \\ &= \text{wt} \left(\sum_{j=1}^{\alpha(\mathcal{H})} \left(\sum_{i \in K} z_i b_{i,j} \right) \mathbf{L}'_j \right) \\ &\geq 2\delta + 1, \end{aligned}$$

where the last transition is due to the existence of $j \in [\alpha(\mathcal{H})]$ such that

$$\sum_{i \in K} z_i b_{i,j} \neq 0,$$

and the fact that \mathbf{L}'_j 's are linearly independent nonzero code-words of a code of minimum distance $2\delta + 1$.

We conclude that the index code based on \mathbf{L} is capable of correcting δ errors. Therefore, $N_q[\mathcal{H}, \delta] = N_q[\alpha(\mathcal{H}), 2\delta + 1]$. \blacksquare

Example 4.4: Let $q = 2$, $m = n = 5$, $f(i) = i$ for all $i \in [m]$, and $\delta = 2$. Assume

$$\begin{aligned} \mathcal{X}_1 &= \{2, 3, 4\}, & \mathcal{X}_2 &= \{3, 4, 5\}, & \mathcal{X}_3 &= \{4, 5, 1\}, \\ \mathcal{X}_4 &= \{5, 1, 2\}, & \mathcal{X}_5 &= \{1, 2, 3\}. \end{aligned}$$

Let $\mathcal{H} = \mathcal{H}(5, 5, \mathcal{X}, f)$. Then

$$\begin{aligned} \mathcal{J}(\mathcal{H}) &= \left\{ \{1\}, \{1, 5\}, \{2\}, \{2, 1\}, \{3\}, \right. \\ &\quad \left. \{3, 2\}, \{4\}, \{4, 3\}, \{5\}, \{5, 4\} \right\}. \end{aligned}$$

It is easy to check that $\alpha(\mathcal{H}) = 2$. Therefore, Theorem 4.3 implies that

$$N_2[\mathcal{H}, 2] \geq N_2[2, 5] = 8.$$

The last equality can be verified by [18].

On the other hand, take the matrix

$$\mathbf{B} \triangleq \begin{pmatrix} 1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 0 & 1 \\ 1 & 1 \end{pmatrix}.$$

The matrix \mathbf{B} satisfies the property that for all $K \in \mathcal{J}(\mathcal{H})$, $K \neq \emptyset$, there exists $j \in [2]$ such that

$$\sum_{i \in K} b_{i,j} \neq 0.$$

From Theorem 4.3, we have $N_2[\mathcal{H}, 2] = N_2[2, 5] = 8$.

Remark 4.5: In [6], when $m = n$ and $f(i) = i$ for all $i \in [n]$, $\alpha(\mathcal{H})$ was shown to be a lower bound on the length of a (non-error-correcting) linear index code. However, the α -bound in Theorem 4.3 does not follow from the results in [6]. The reason is that a concatenation of an optimal linear error-correcting code with an optimal non-error-correcting index code might fail to produce an optimal linear ECIC. This is illustrated later in Example 4.8.

The following proposition is based on the fact that concatenation of a δ -error-correcting code with an optimal (non-error-correcting) \mathcal{H} -IC yields a (δ, \mathcal{H}) -ECIC.

Proposition 4.6 (κ -bound): The length of an optimal (δ, \mathcal{H}) -ECIC over \mathbb{F}_q satisfies

$$\mathcal{N}_q[\mathcal{H}, \delta] \leq N_q[\kappa_q(\mathcal{H}), 2\delta + 1].$$

Proof: Let \mathbf{G} , which is an $n \times \kappa_q(\mathcal{H})$ matrix, correspond to an optimal \mathcal{H} -IC over \mathbb{F}_q . Denote

$$\mathbf{y} = \mathbf{x}\mathbf{G} \in \mathbb{F}_q^{\kappa_q(\mathcal{H})}.$$

Let \mathbf{M} be a generator matrix of an optimal $[N, \kappa_q(\mathcal{H}), 2\delta + 1]_q$ code \mathcal{C}' , where

$$N = N_q[\kappa_q(\mathcal{H}), 2\delta + 1].$$

Consider a scheme where S broadcasts the vector $\mathbf{y}\mathbf{M} \in \mathbb{F}_q^N$. If less than δ errors occur, then each receiver R_i is able to recover \mathbf{y} by using \mathcal{C}' . Hence each R_i is able to recover $x_{f(i)}$. Therefore, for the index code based on \mathbf{L} ,

$$\mathbf{L} = \mathbf{G}\mathbf{M},$$

each receiver R_i is capable to recover $x_{f(i)}$ if the number of errors is less or equal to δ . The length of the corresponding ECIC is $N = N_q[\kappa_q(\mathcal{H}), 2\delta + 1]$. Therefore,

$$\mathcal{N}_q[\mathcal{H}, \delta] \leq N_q[\kappa_q(\mathcal{H}), 2\delta + 1].$$

■

By combining the results in Theorem 4.3 and in Proposition 4.6, we obtain the following corollary.

Corollary 4.7: The length of an optimal linear (δ, \mathcal{H}) -ECIC over \mathbb{F}_q satisfies

$$N_q[\alpha(\mathcal{H}), 2\delta + 1] \leq \mathcal{N}_q[\mathcal{H}, \delta] \leq N_q[\kappa_q(\mathcal{H}), 2\delta + 1].$$

It is shown in the example below that the inequalities in Corollary 4.7 can be strict. In particular, it follows that mere application of an error-correcting code on top of an index code may fail to provide us with an optimal linear ECIC. This fact motivates the study of ECIC's in Sections III–VII.

Example 4.8: Let $q = 2$, $m = n = 5$, $\delta = 2$, and $f(i) = i$ for all $i \in [m]$. Assume

$$\begin{aligned} \mathcal{X}_1 &= \{2, 5\}, & \mathcal{X}_2 &= \{1, 3\}, & \mathcal{X}_3 &= \{2, 4\}, \\ & & \mathcal{X}_4 &= \{3, 5\}, & \mathcal{X}_5 &= \{1, 4\}. \end{aligned}$$

Let $\mathcal{H} = \mathcal{H}(5, 5, \mathcal{X}, f)$. Then we have

$$\begin{aligned} \mathcal{J}(\mathcal{H}) = \{ & \{1\}, \{1, 3\}, \{1, 4\}, \{1, 3, 4\}, \\ & \{2\}, \{2, 4\}, \{2, 5\}, \{2, 4, 5\}, \\ & \{3\}, \{1, 3\}, \{3, 5\}, \{1, 3, 5\}, \\ & \{4\}, \{1, 4\}, \{2, 4\}, \{1, 2, 4\}, \\ & \{5\}, \{2, 5\}, \{3, 5\}, \{2, 3, 5\} \}. \end{aligned}$$

The side information graph $\mathcal{G}_{\mathcal{H}}$ of this instance is a pentagon. It is easy to verify that $\alpha(\mathcal{H}) = \alpha(\mathcal{G}_{\mathcal{H}}) = 2$. It follows from

Theorem 9 in [7] that $\kappa_2(\mathcal{H}) = \min\text{-rank}_2(\mathcal{G}_{\mathcal{H}}) = 3$. Thus, from [18] we have

$$N_2[2, 5] = 8 \quad \text{and} \quad N_2[3, 5] = 10.$$

Due to Corollary 4.7, we have

$$8 \leq \mathcal{N}_2[\mathcal{H}, 2] \leq 10.$$

Using a computer search, we obtain that $\mathcal{N}_2[\mathcal{H}, 2] = 9$, and the corresponding optimal scheme is based on

$$\mathbf{L} = \begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 1 & 1 & 0 & 0 & 1 & 0 & 1 & 1 \\ 1 & 0 & 1 & 0 & 1 & 0 & 0 & 1 & 1 \end{pmatrix}.$$

It is technical to verify that for all $K \in \mathcal{J}(\mathcal{H})$,

$$\text{wt}\left(\sum_{i \in K} \mathbf{L}_i\right) \geq 5.$$

Therefore by Lemma 3.8, for the index code based on \mathbf{L} , each receiver R_i is able to recover x_i , if the number of errors is less than or equal to 2. Observe that the length of the ECIC corresponding to \mathbf{L} lies strictly between the α -bound and the κ -bound.

When the graph \mathcal{G} is undirected (or symmetric), the following theorem holds (see, for instance, [16]).

Theorem 4.9: Let $\chi(\bar{\mathcal{G}})$ denote the chromatic number of the complement of the graph \mathcal{G} . Then,

$$\alpha(\mathcal{G}) \leq \min\text{-rank}_q(\mathcal{G}) \leq \chi(\bar{\mathcal{G}}).$$

When $m = n$ and $f(i) = i$ for all $i \in [m]$, we have that $\alpha(\mathcal{H}) = \alpha(\mathcal{G}_{\mathcal{H}})$ and $\kappa_q(\mathcal{H}) = \min\text{-rank}_q(\mathcal{G}_{\mathcal{H}})$. Moreover, if the graph $\mathcal{G}_{\mathcal{H}}$ is symmetric and satisfies $\alpha(\mathcal{G}_{\mathcal{H}}) = \chi(\bar{\mathcal{G}}_{\mathcal{H}})$, then from Corollary 4.7 we have

$$\mathcal{N}_q[\mathcal{H}, \delta] = N_q[\alpha(\mathcal{H}), 2\delta + 1] = N_q[\kappa_q(\mathcal{H}), 2\delta + 1],$$

for all q , and the corresponding bounds in Corollary 4.7 are tight.

Definition 4.10: An undirected (or symmetric) graph \mathcal{G} is called perfect if for every induced subgraph \mathcal{G}' of \mathcal{G} , $\alpha(\mathcal{G}') = \chi(\bar{\mathcal{G}}')$.

Perfect graphs include families of graphs such as trees, bipartite graphs, interval graphs, and chordal graphs. If $m = n$, $f(i) = i$ for all $i \in [m]$, and $\mathcal{G}_{\mathcal{H}}$ is perfect, then the bounds in Corollary 4.7 are tight. For the full characterization of perfect graphs, the reader can refer to [19].

V. THE SINGLETON BOUND

The following bound is analogous to Singleton bound for classical linear error-correcting codes.

Theorem 5.1 (Singleton bound): The length of an optimal linear (δ, \mathcal{H}) -ECIC over \mathbb{F}_q satisfies

$$\mathcal{N}_q[\mathcal{H}, \delta] \geq \kappa_q(\mathcal{H}) + 2\delta.$$

Proof: Let \mathbf{L} be the $n \times \mathcal{N}_q[\mathcal{H}, \delta]$ matrix corresponding to some optimal (δ, \mathcal{H}) -ECIC. Let \mathbf{L}' be the matrix obtained by deleting any 2δ columns from \mathbf{L} .

By Lemma 3.8, \mathbf{L} satisfies

$$\text{wt} \left(\sum_{i \in K} z_i \mathbf{L}_i \right) \geq 2\delta + 1 ,$$

for all $K \in \mathcal{J}(\mathcal{H})$ and all choices of $z_i \in \mathbb{F}_q^*$, $i \in K$. We deduce that the rows of \mathbf{L}' also satisfy

$$\text{wt} \left(\sum_{i \in K} z_i \mathbf{L}'_i \right) \geq 1 .$$

By Corollary 3.10, \mathbf{L}' corresponds to a linear \mathcal{H} -IC. Therefore, by Lemma 3.5, part 2, \mathbf{L}' has at least $\kappa_q(\mathcal{H})$ columns. We deduce that

$$\mathcal{N}_q[\mathcal{H}, \delta] - 2\delta \geq \kappa_q(\mathcal{H}) ,$$

which concludes the proof. \blacksquare

The following corollary from Proposition 4.6 and Theorem 5.1 demonstrates that, for sufficiently large alphabets, a concatenation of a classical MDS error-correcting code with an optimal (non-error-correcting) index code yields an optimal ECIC. However, as it was illustrated in Example 4.8, this does not hold for the index coding schemes over small alphabets.

Corollary 5.2 (MDS error-correcting index code): For $q \geq \kappa_q(\mathcal{H}) + 2\delta - 1$,

$$\mathcal{N}_q[\mathcal{H}, \delta] = \kappa_q(\mathcal{H}) + 2\delta . \quad (8)$$

Proof: From Theorem 5.1, we have

$$\mathcal{N}_q[\mathcal{H}, \delta] \geq \kappa_q(\mathcal{H}) + 2\delta .$$

On the other hand, from Proposition 4.6,

$$\mathcal{N}_q[\mathcal{H}, \delta] \leq \mathcal{N}_q[\kappa_q(\mathcal{H}), 2\delta + 1] = \kappa_q(\mathcal{H}) + 2\delta ,$$

for $q \geq \kappa_q(\mathcal{H}) + 2\delta - 1$ (by taking doubly-extended Reed-Solomon (RS) codes). Therefore, for these q , (8) holds. \blacksquare

Remark 5.3: Let $q = 2$, $m = n = 2\ell + 1$ ($\ell \geq 2$), and $f(i) = i$ for all $i \in [m]$. Let $\mathcal{X}_1 = \{2, n\}$ and $\mathcal{X}_n = \{1, n-1\}$. For $2 \leq i \leq n$, let $\mathcal{X}_i = \{i-1, i+1\}$. Let $\mathcal{H} = \mathcal{H}(n, n, \mathcal{X}, f)$. Then $\mathcal{G}_{\mathcal{H}}$ is the (symmetric) odd cycle of length n . Therefore, $\alpha(\mathcal{H}) = \alpha(\mathcal{G}_{\mathcal{H}}) = \ell$. From [7], $\kappa_2(\mathcal{H}) = \min\text{-rank}_2(\mathcal{G}_{\mathcal{H}}) = \ell + 1$. From α -bound,

$$\mathcal{N}_2[\mathcal{H}, \delta] \geq \mathcal{N}_2[\ell, 2\delta + 1] .$$

By contrast, from Theorem 5.1,

$$\mathcal{N}_2[\mathcal{H}, \delta] \geq (\ell + 1) + 2\delta .$$

As there are no nontrivial binary MDS codes, we have

$$\mathcal{N}_2[\ell, 2\delta + 1] \geq \ell + 2\delta + 1 ,$$

for all choices of $\delta > 0$. Therefore, for these choices, the α -bound is at least as good as the Singleton bound.

VI. RANDOM CODES

In this section we prove an inexplicit upper bound on the optimal length of the ECIC's. The proof is based on constructing a random ECIC and analyzing its parameters.

Theorem 6.1: Let $\mathcal{H} = \mathcal{H}(m, n, \mathcal{X}, f)$ describe an instance of the ICSI problem. Then there exists a (δ, \mathcal{H}) -ECIC over \mathbb{F}_q of length N if

$$\sum_{i \in [m]} q^{n - |\mathcal{X}_i| - 1} < \frac{q^N}{V_q(N, 2\delta)} ,$$

where

$$V_q(N, 2\delta) = \sum_{\ell=0}^{2\delta} \binom{N}{\ell} (q-1)^\ell \quad (9)$$

is the volume of the q -ary sphere in \mathbb{F}_q^N .

Proof: We construct a random $n \times N$ matrix \mathbf{L} over \mathbb{F}_q , row by row. Each row is selected independently of other rows, uniformly over \mathbb{F}_q^N . Define vector spaces

$$\mathbf{M}_i \triangleq \text{span}(\{\mathbf{L}_j : j \in \mathcal{Y}_i\})$$

for all $i \in [m]$. We also define the following events:

$$\forall i \in [m] : \text{Event } E_i \triangleq \{d(\mathbf{L}_{f(i)}, \mathbf{M}_i) < 2\delta + 1\} ,$$

and

$$\text{Event } E_{\text{Fail}} \triangleq$$

$$\{\mathbf{L} \text{ does not correspond to a } (\delta, \mathcal{H})\text{-ECIC}\} .$$

The event E_i represents the situation when the receiver R_i cannot recover $x_{f(i)}$. Then, by Corollary 3.9, the event E_{Fail} is equivalent to $\bigcup_{i \in [m]} E_i$. Therefore,

$$\Pr(E_{\text{Fail}}) = \Pr\left(\bigcup_{i \in [m]} E_i\right) \leq \sum_{i \in [m]} \Pr(E_i) . \quad (10)$$

For a particular event E_i , $i \in [m]$,

$$\Pr(E_i) \leq \frac{q^{|\mathcal{Y}_i|} V_q(N, 2\delta)}{q^N} . \quad (11)$$

There exists a matrix \mathbf{L} that corresponds to a (δ, \mathcal{H}) -ECIC if $\Pr(E_{\text{Fail}}) < 1$. It is enough to require that the right-hand side of (10) is smaller than 1. By plugging in the expression in (11), we obtain a sufficient condition on the existence of a (δ, \mathcal{H}) -ECIC over \mathbb{F}_q :

$$\frac{V_q(N, 2\delta)}{q^N} \sum_{i \in [m]} q^{|\mathcal{Y}_i|} < 1 .$$

\blacksquare

Remark 6.2: The bound in Theorem 6.1 does not take into account the structure of the sets \mathcal{X}_i 's, other than their cardinalities. Therefore, this bound generally is weaker than the κ -bound. On the other hand, for a particular instance of the ICSI problem, it is easier to compute this bound, while calculating the κ -bound in general is an NP-hard problem.

Remark 6.3: The bound in Theorem 6.1 implies a bound on $\kappa_q(\mathcal{H})$, which is tight for some \mathcal{H} . Indeed, fix $\delta = 0$. The bound implies that there exists a linear index code of length N whenever

$$\sum_{i \in [m]} q^{n-|\mathcal{X}_i|-1} < q^N. \quad (12)$$

Let $m = n = 2\ell + 1$ ($\ell \geq 2$), and $f(i) = i$ for all $i \in [n]$. Let $\mathcal{X}_1 = [n] \setminus \{1, 2, n\}$ and $\mathcal{X}_n = [n] \setminus \{1, n-1, n\}$. For $2 \leq i \leq n-1$, let $\mathcal{X}_i = [n] \setminus \{i-1, i, i+1\}$. Let $\mathcal{H} = \mathcal{H}(n, n, \mathcal{X}, f)$ be the corresponding side information hypergraph. Then $\mathcal{G}_{\mathcal{H}}$ is the complement of the (symmetric directed) odd cycle of length n . We have $|\mathcal{X}_i| = 2\ell - 2$ for all $i \in [n]$. Then (12) becomes

$$N > 2 + \log_q(2\ell + 1).$$

If $q > 2\ell + 1$ then we obtain $N \geq 3$. Observe that in this case $\kappa_q(\mathcal{H}) = \min\text{-rank}_q(\mathcal{G}_{\mathcal{H}}) = 3$ (see [13, Claim A.1]), and thus the bound is tight.

VII. SYNDROME DECODING

Consider the (δ, \mathcal{H}) -ECIC based on a matrix \mathbf{L} . Suppose that the receiver R_i , $i \in [m]$, receives the vector

$$\mathbf{y}_i = \mathbf{x}\mathbf{L} + \boldsymbol{\epsilon}_i, \quad (13)$$

where $\mathbf{x}\mathbf{L}$ is the codeword transmitted by S , and $\boldsymbol{\epsilon}_i$ is the error pattern affecting this codeword.

In the classical coding theory, the transmitted vector \mathbf{c} , the received vector \mathbf{y} , and the error pattern \mathbf{e} are related by $\mathbf{y} = \mathbf{c} + \mathbf{e}$. Therefore, if \mathbf{y} is known to the receiver, then there is a one-to-one correspondence between the values of unknown vectors \mathbf{c} and \mathbf{e} . For index coding, however, this is no longer the case. The following theorem shows that, in order to recover the message $x_{f(i)}$ from \mathbf{y}_i using (13), it is sufficient to find just one vector from a set of possible error patterns. This set is defined as follows:

$$\mathcal{L}_i(\boldsymbol{\epsilon}_i) = \{\boldsymbol{\epsilon}_i + \mathbf{z} : \mathbf{z} \in \text{span}(\{\mathbf{L}_j\}_{j \in \mathcal{Y}_i})\}.$$

We henceforth refer to the set $\mathcal{L}_i(\boldsymbol{\epsilon}_i)$ as the *set of relevant error patterns*.

Lemma 7.1: Assume that the receiver R_i receives \mathbf{y}_i .

- 1) If R_i knows the message $x_{f(i)}$ then it is able to determine the set $\mathcal{L}_i(\boldsymbol{\epsilon}_i)$.
- 2) If R_i knows some vector $\hat{\mathbf{e}} \in \mathcal{L}_i(\boldsymbol{\epsilon}_i)$ then it is able to determine $x_{f(i)}$.

Proof:

- 1) From (13), we have

$$\mathbf{y}_i = x_{f(i)}\mathbf{L}_{f(i)} + \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i} + \mathbf{x}_{\mathcal{Y}_i}\mathbf{L}_{\mathcal{Y}_i} + \boldsymbol{\epsilon}_i. \quad (14)$$

If R_i knows $x_{f(i)}$, then it is also able to determine

$$\boldsymbol{\epsilon}_i + \mathbf{x}_{\mathcal{Y}_i}\mathbf{L}_{\mathcal{Y}_i} = \mathbf{y}_i - x_{f(i)}\mathbf{L}_{f(i)} - \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i} \in \mathcal{L}_i(\boldsymbol{\epsilon}_i).$$

Since R_i has a knowledge of \mathbf{L} , it is also able to determine the whole $\mathcal{L}_i(\boldsymbol{\epsilon}_i)$.

- 2) Suppose that R_i knows a vector

$$\hat{\mathbf{e}} = \boldsymbol{\epsilon}_i + \sum_{j \in \mathcal{Y}_i} z_j \mathbf{L}_j \in \mathcal{L}_i(\boldsymbol{\epsilon}_i),$$

for some $\mathbf{z} = (z_j)_{j \in \mathcal{Y}_i} \in \mathbb{F}_q^{|\mathcal{Y}_i|}$. We show that R_i is able then to determine $x_{f(i)}$. Indeed, we re-write (14) as

$$\mathbf{y}_i = x_{f(i)}\mathbf{L}_{f(i)} + \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i} + (\mathbf{x}_{\mathcal{Y}_i} - \mathbf{z})\mathbf{L}_{\mathcal{Y}_i} + \hat{\mathbf{e}}. \quad (15)$$

The receiver R_i can find some solution of the equation

$$\mathbf{y}_i = \hat{x}_{f(i)}\mathbf{L}_{f(i)} + \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i} + \hat{\mathbf{x}}_{\mathcal{Y}_i}\mathbf{L}_{\mathcal{Y}_i} + \hat{\mathbf{e}}, \quad (16)$$

with respect to the unknowns $\hat{x}_{f(i)}$ and $\hat{\mathbf{x}}_{\mathcal{Y}_i}$. Observe that (16) has at least one solution due to (15).

From (15) and (16), we deduce that

$$\mathbf{0} = (\hat{x}_{f(i)} - x_{f(i)})\mathbf{L}_{f(i)} + (\hat{\mathbf{x}}_{\mathcal{Y}_i} - \mathbf{x}_{\mathcal{Y}_i} + \mathbf{z})\mathbf{L}_{\mathcal{Y}_i}.$$

This equality implies that $\hat{x}_{f(i)} = x_{f(i)}$ (otherwise, by Corollary 3.9, the sum in the right-hand side will have nonzero weight). Hence, R_i is able to determine $x_{f(i)}$, as claimed. ■

We now describe a syndrome decoding algorithm for linear error-correcting index codes. From (14), we have

$$\mathbf{y}_i - \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i} - \boldsymbol{\epsilon}_i \in \text{span}(\{\mathbf{L}_{f(i)}\} \cup \{\mathbf{L}_j\}_{j \in \mathcal{Y}_i}).$$

Let $\mathcal{C}_i = \text{span}(\{\mathbf{L}_{f(i)}\} \cup \{\mathbf{L}_j\}_{j \in \mathcal{Y}_i})$, and let $\mathbf{H}^{(i)}$ be a parity check matrix of \mathcal{C}_i . We obtain that

$$\mathbf{H}^{(i)}\boldsymbol{\epsilon}_i^T = \mathbf{H}^{(i)}(\mathbf{y}_i - \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i})^T. \quad (17)$$

Let $\boldsymbol{\beta}_i$ be a column vector defined by

$$\boldsymbol{\beta}_i = \mathbf{H}^{(i)}(\mathbf{y}_i - \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i})^T. \quad (18)$$

Observe that each R_i is capable of determining $\boldsymbol{\beta}_i$. Then we can re-write (17) as

$$\mathbf{H}^{(i)}\boldsymbol{\epsilon}_i^T = \boldsymbol{\beta}_i.$$

This leads us to the formulation of the following decoding procedure for R_i .

• *Input:* $\mathbf{y}_i, \mathbf{x}_{\mathcal{X}_i}, \mathbf{L}$.

• *Step 1:* Compute the syndrome

$$\boldsymbol{\beta}_i = \mathbf{H}^{(i)}(\mathbf{y}_i - \mathbf{x}_{\mathcal{X}_i}\mathbf{L}_{\mathcal{X}_i})^T.$$

• *Step 2:* Find the lowest Hamming weight solution $\hat{\mathbf{e}}$ of the system

$$\mathbf{H}^{(i)}\hat{\mathbf{e}}^T = \boldsymbol{\beta}_i. \quad (19)$$

• *Step 3:* Given that $\hat{\mathbf{x}}_{\mathcal{X}_i} = \mathbf{x}_{\mathcal{X}_i}$, solve the system for $\hat{x}_{f(i)}$:

$$\mathbf{y}_i = \hat{\mathbf{x}}\mathbf{L} + \hat{\mathbf{e}}. \quad (20)$$

• *Output:* $\hat{x}_{f(i)}$.

Fig. 2: Syndrome decoding procedure.

Remark 7.2: Gaussian elimination can be used to solve (20) for $\hat{x}_{f(i)}$. However, since \mathbf{L} also corresponds to an \mathcal{H} -IC, there is more efficient way to do so. From Lemma 3.5, there exists

a vector $\mathbf{v}_i \triangleleft \mathcal{X}_i$ satisfying $\mathbf{v}_i + \mathbf{e}_{f(i)} \in \text{colspan}(\mathbf{L})$. Hence $\mathbf{v}_i + \mathbf{e}_{f(i)} = \mathbf{u}\mathbf{L}^T$ for some $\mathbf{u} \in \mathbb{F}_q^N$. Therefore

$$\begin{aligned}\hat{\mathbf{x}}_{f(i)} &= \hat{\mathbf{x}}(\mathbf{v}_i + \mathbf{e}_{f(i)})^T - \hat{\mathbf{x}}\mathbf{v}_i^T \\ &= \hat{\mathbf{x}}\mathbf{L}\mathbf{u}^T - \hat{\mathbf{x}}\mathbf{v}_i^T \\ &= (\mathbf{y}_i - \hat{\boldsymbol{\epsilon}})\mathbf{u}^T - \hat{\mathbf{x}}\mathbf{v}_i^T.\end{aligned}$$

With the knowledge of \mathbf{L} and $\mathbf{x}_{\mathcal{X}_i}$, R_i can determine \mathbf{u} and $\hat{\mathbf{x}}\mathbf{v}_i^T$. Therefore, it can also determine $\hat{\mathbf{x}}_{f(i)}$. Note that (20) may have more than one solution $\hat{\mathbf{x}}$ with $\hat{\mathbf{x}}_{\mathcal{X}_i} = \mathbf{x}_{\mathcal{X}_i}$. However, as shown in the next theorem, if at most δ errors occur in \mathbf{y}_i , then it always holds that $\hat{\mathbf{x}}_{f(i)} = x_{f(i)}$.

Theorem 7.3: Let $\mathbf{y}_i = \mathbf{x}\mathbf{L} + \boldsymbol{\epsilon}_i$ be the vector received by R_i , and let $\text{wt}(\boldsymbol{\epsilon}_i) \leq \delta$. Assume that the procedure in Figure 2 is applied to $(\mathbf{y}_i, \mathbf{x}_{\mathcal{X}_i}, \mathbf{L})$. Then, its output satisfies $\hat{\mathbf{x}}_{f(i)} = x_{f(i)}$.

Proof: By Lemma 7.1, it is sufficient to prove that $\hat{\boldsymbol{\epsilon}} \in \mathcal{L}_i(\boldsymbol{\epsilon}_i)$. Indeed, since

$$\mathbf{H}^{(i)}\boldsymbol{\epsilon}_i^T = \mathbf{H}^{(i)}\hat{\boldsymbol{\epsilon}}^T = \boldsymbol{\beta}_i,$$

we have

$$\mathbf{H}^{(i)}(\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}_i)^T = \mathbf{0}.$$

Hence, $\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}_i \in \mathcal{C}_i$, and therefore,

$$\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}_i = z_{f(i)}\mathbf{L}_{f(i)} + \sum_{j \in \mathcal{Y}_i} z_j \mathbf{L}_j, \quad (21)$$

for some $z_{f(i)} \in \mathbb{F}_q$ and $z_j \in \mathbb{F}_q$, $j \in \mathcal{Y}_i$.

Since $\boldsymbol{\epsilon}_i$ is a solution of (19), and $\text{wt}(\boldsymbol{\epsilon}_i) \leq \delta$, we deduce that $\text{wt}(\hat{\boldsymbol{\epsilon}}) \leq \delta$ as well. Hence,

$$\text{wt}\left(z_{f(i)}\mathbf{L}_{f(i)} + \sum_{j \in \mathcal{Y}_i} z_j \mathbf{L}_j\right) = \text{wt}(\hat{\boldsymbol{\epsilon}} - \boldsymbol{\epsilon}_i) \leq 2\delta.$$

Therefore, by Corollary 3.9, $z_{f(i)} = 0$. Hence, $\hat{\boldsymbol{\epsilon}} \in \mathcal{L}_i(\boldsymbol{\epsilon}_i)$, as desired, and therefore $\hat{\mathbf{x}}_{f(i)} = x_{f(i)}$. ■

Remark 7.4: We anticipate Step 2 in Figure 2 to be computationally hard. Indeed, the problem of finding $\hat{\boldsymbol{\epsilon}}$ over \mathbb{F}_2 of the lowest weight satisfying

$$\mathbf{H}^{(i)}\hat{\boldsymbol{\epsilon}}^T = \boldsymbol{\beta}_i, \quad (22)$$

for a given binary vector $\boldsymbol{\beta}_i$ is at least as hard as a decision problem COSET WEIGHTS that was shown in [20] to be NP-complete.

VIII. STATIC CODES AND RELATED PROBLEMS

A. Static Error-Correcting Index Codes

In the previous sections we focused on linear δ -error-correcting index codes for a *particular* instance of the ICSI problem. When some of the parameters m , n , \mathcal{X} , and f are variable or not known, it is very likely that an error-correcting index code for the instance with particular values of these parameters can not be used for the instances with different values of some of these parameters. Therefore, it is interesting to design an error-correcting index code which will be suitable for a *family* of instances of the ICSI problem.

Definition 8.1: Let $\Gamma = \{(m, n, \mathcal{X}, f)\}$ be a set of instances for an ICSI problem. A δ -error-correcting index code over \mathbb{F}_q is said to be *static* under the set Γ if it is a δ -error-correcting (m, n, \mathcal{X}, f) -IC over \mathbb{F}_q for all instances $(m, n, \mathcal{X}, f) \in \Gamma$.

Recall that an instance (m, n, \mathcal{X}, f) can be described by the side information hypergraph $\mathcal{H}(m, n, \mathcal{X}, f)$. For a set Γ of instances (m, n, \mathcal{X}, f) , let

$$\mathfrak{J}(\Gamma) \triangleq \bigcup_{(m, n, \mathcal{X}, f) \in \Gamma} \mathcal{J}(\mathcal{H}(m, n, \mathcal{X}, f)), \quad (23)$$

where $\mathcal{J}(\mathcal{H}(m, n, \mathcal{X}, f))$ is defined as in (3). We also define

$$n(\Gamma) \triangleq \max\{n : (m, n, \mathcal{X}, f) \in \Gamma\}.$$

Lemma 8.2: The $n(\Gamma) \times N$ matrix \mathbf{L} corresponds to a δ -error-correcting index code which is static under Γ if and only if

$$\text{wt}\left(\sum_{i \in K} z_i \mathbf{L}_i\right) \geq 2\delta + 1,$$

for all $K \in \mathfrak{J}(\Gamma)$ and for all choices of $z_i \in \mathbb{F}_q^*$, $i \in K$.

Proof: The proof follows from Definition 8.1 and Lemma 3.8. ■

Please notice that when \mathbf{L} is used for an instance $(m, n, \mathcal{X}, f) \in \Gamma$ with $n < n(\Gamma)$, then the last $n(\Gamma) - n$ rows of \mathbf{L} are simply discarded.

One particular family of interest is $\Gamma(n, \rho)$, the family that contains all instances where each receiver owns at least $n - \rho$ messages as its side information. More formally,

$$\begin{aligned}\Gamma(n, \rho) &= \{(m, n', \mathcal{X}, f) : n' \leq n \\ &\quad \text{and } \forall i \in [m], |\mathcal{X}_i| \geq n - \rho\}.\end{aligned}$$

A δ -error-correcting index code which is static under $\Gamma(n, \rho)$ will provide successful communication between the sender and the receivers under the presence of at most δ errors, despite a possible change of the collection of the side information sets \mathcal{X} , a change of the set of receivers, and a change of the demand function, as long as each receiver still possesses at least $n - \rho$ messages.

In the rest of this section, we assume that $N \geq 1$, $n \geq \rho \geq 1$ and $\delta \geq 0$.

Definition 8.3: An $n \times N$ matrix \mathbf{L} is said to satisfy the (ρ, δ) -Property if any nontrivial linear combination of at most ρ rows of \mathbf{L} has weight at least $2\delta + 1$.

Proposition 8.4: The $n \times N$ matrix \mathbf{L} corresponds to a δ -error-correcting linear index code, which is static under $\Gamma(n, \rho)$, if and only if \mathbf{L} satisfies the (ρ, δ) -Property.

Proof: Let \mathbf{L} be an $n \times N$ matrix that satisfies the (ρ, δ) -Property. We show that this is equivalent to the condition that \mathbf{L} corresponds to a δ -error-correcting linear index code, which is static under $\Gamma(n, \rho)$. By Lemma 8.2, it suffices to show that $\mathfrak{J}(\Gamma(n, \rho))$ is the collection of all nonempty subsets of $[n]$, whose cardinalities are not greater than ρ .

Consider an instance $(m, n', \mathcal{X}, f) \in \Gamma(n, \rho)$. For all $i \in [m]$, we have $|\mathcal{X}_i| \geq n - \rho$ and $\mathcal{Y}_i = [n'] \setminus (f(i) \cup \mathcal{X}_i)$, and thus we deduce that

$$|\mathcal{Y}_i| \leq n' - 1 - (n - \rho) \leq n' - 1 - (n' - \rho) = \rho - 1.$$

Hence by (3), the cardinality of each set in $\mathcal{J}(\mathcal{H}(m, n', \mathcal{X}, f))$ is at most

$$1 + (\rho - 1) = \rho.$$

Therefore, due to (23), every set in $\mathfrak{J}(\Gamma(n, \rho))$ has at most ρ elements.

It remains to show that every nonempty subset of $[n]$ whose cardinality is at most ρ belongs to $\mathfrak{J}(\Gamma(n, \rho))$. Consider an arbitrary ρ' -subset $K = \{i_1, i_2, \dots, i_{\rho'}\}$ of $[n]$, with $1 \leq \rho' \leq \rho$. Consider an instance $(m = 1, n, \mathcal{X}, f) \in \Gamma(n, \rho)$ with $\mathcal{X}_1 = [n] \setminus K$ and $f(1) = i_1$. Since

$$\mathcal{Y}_1 = K \setminus \{i_1\},$$

we have

$$K = \{i_1\} \cup \mathcal{Y}_1 \in \mathcal{J}(\mathcal{H}(m, n, \mathcal{X}, f)) \subseteq \mathfrak{J}(\Gamma(n, \rho)).$$

The proof follows. \blacksquare

B. Application: Weakly Resilient Functions

In this section we introduce the notion of weakly resilient functions. Hereafter, we restrict the discussion to the binary alphabet.

The concept of *binary resilient functions* was first introduced by Chor *et. al.* in [21] and independently by Bennet *et. al.* in [22].

Definition 8.5: A function $\mathbf{f} : \mathbb{F}_2^N \rightarrow \mathbb{F}_2^n$ is called *t-resilient* if \mathbf{f} satisfies the following property: when t arbitrary inputs of \mathbf{f} are fixed and the remaining $N - t$ inputs run through all the 2^{N-t} -tuples exactly once, the value of \mathbf{f} runs through every possible output n -tuple an equal number of times. Moreover, if \mathbf{f} is a linear transformation then it is called a *linear t-resilient function*. We refer to the parameter t as the *resiliency* of \mathbf{f} .

The applications of resilient functions can be found in fault-tolerant distributed computing, quantum cryptographic key distribution [21], privacy amplification [22] and random sequence generation for stream ciphers [23]. Connections between linear error-correcting codes and resilient functions were established in [21].

Theorem 8.6 ([21]): Let \mathbf{L} be an $n \times N$ binary matrix. Then \mathbf{L} is a generator matrix of a linear error-correcting code with minimum distance $d = t + 1$ if and only if $\mathbf{f}(z) = \mathbf{L}z^T$ is *t-resilient*.

Remark 8.7: Vectorial boolean functions with certain properties are useful for design of stream ciphers. These properties include high resiliency and high nonlinearity (see, for instance, [23]). However, linear resilient functions are still particularly interesting, since they can be transformed into highly nonlinear resilient functions with the same parameters. This can be

achieved by a composition of the linear function with a highly nonlinear permutation (see [24], [25] for more details).

Below we introduce a definition of a ρ -weakly *t-resilient* function, which is a weaker version of a *t-resilient* function.

Definition 8.8: A function $\mathbf{f} : \mathbb{F}_2^N \rightarrow \mathbb{F}_2^n$ is called *ρ -weakly t-resilient* if \mathbf{f} satisfies the property that every set of ρ coordinates in the image of \mathbf{f} runs through every possible output ρ -tuple an equal number of times, when t arbitrary inputs of \mathbf{f} are fixed and the remaining $N - t$ inputs run through all the 2^{N-t} -tuples exactly once.

Remark 8.9: A ρ -weakly *t-resilient* function $\mathbf{f} : \mathbb{F}_2^N \rightarrow \mathbb{F}_2^n$ can be viewed as a collection of $\binom{n}{\rho}$ different *t-resilient* functions $\mathbb{F}_2^N \rightarrow \mathbb{F}_2^\rho$, each such function is obtained by taking some ρ coordinates in the image of \mathbf{f} . Similarly to [21], consider a scenario, in which two parties are sharing a secret key, which consists of N randomly selected bits. Suppose that at some moment t out of the N bits of the key are leaked to an adversary. By applying a *t-resilient* function to the current N -bit key, two parties are able to obtain a completely new and secret key of n bits, without requiring any communication or randomness generation. However, if the parties use various parts of the key for various purposes, they may only require one of the ρ -bit secret keys (instead of the larger n -bit key). In that case a ρ -weakly *t-resilient* function can be used. By applying a ρ -weakly *t-resilient* function to the current N -bit key, the parties obtain a set of $\binom{n}{\rho}$ different ρ -bit keys, each key is new and secret (however these keys might not be independent of each other).

Theorem 8.10: Let \mathbf{L} be an $n \times N$ binary matrix. Then \mathbf{L} satisfies the (ρ, δ) -Property if and only if the function $\mathbf{f} : \mathbb{F}_2^N \rightarrow \mathbb{F}_2^n$ defined by $\mathbf{f}(z) = \mathbf{L}z^T$ is ρ -weakly 2δ -resilient.

Proof:

- 1) Suppose that \mathbf{L} satisfies the (ρ, δ) -Property. Take any ρ -subset $K \subseteq [n]$. By Definition 8.3, the $\rho \times N$ submatrix \mathbf{L}_K of \mathbf{L} is a generating matrix of the error-correcting code with the minimum distance $\geq 2\delta + 1$. By Theorem 8.6, the function $\mathbf{f}_K : \mathbb{F}_2^N \rightarrow \mathbb{F}_2^\rho$ defined by $\mathbf{f}_K(z) = \mathbf{L}_K z^T$ is 2δ -resilient. Since K is an arbitrary ρ -subset of $[n]$, the function \mathbf{f} is ρ -weakly 2δ -resilient.
- 2) Conversely, assume that the function \mathbf{f} is ρ -weakly 2δ -resilient. Take any subset $K \subseteq [n]$, $|K| = \rho$. Then the function $\mathbf{f}_K : \mathbb{F}_2^N \rightarrow \mathbb{F}_2^\rho$ defined by $\mathbf{f}_K(z) = \mathbf{L}_K z^T$ is 2δ -resilient. Therefore, by Theorem 8.6, \mathbf{L}_K is a generating matrix of a linear code with minimum distance $2\delta + 1$. Since K is an arbitrary ρ -subset of $[n]$, by Proposition 8.4 \mathbf{L} satisfies the (ρ, δ) -Property. \blacksquare

C. Bounds and Constructions

In this section we study the problem of constructing a matrix \mathbf{L} satisfying the (ρ, δ) -Property. Such \mathbf{L} with the minimal possible number of columns is called *optimal*. First, observe

that from Proposition 8.4 we have

$$\mathfrak{J}(\Gamma(n, \rho)) = \bigcup_{i=1}^{\rho} \binom{[n]}{i},$$

is the set of all nonempty subsets of $[n]$ of cardinality at most ρ . Next, consider an instance $(m^*, n, \mathcal{X}^*, f^*)$ satisfying

$$\mathcal{J}(\mathcal{H}^*) = \mathfrak{J}(\Gamma(n, \rho)), \quad (24)$$

where $\mathcal{H}^* = \mathcal{H}(m^*, n, \mathcal{X}^*, f^*)$ is the side information hypergraph corresponding to that instance. Such an instance can be constructed as follows. For each subset $K = \{i_1, i_2, \dots, i_{\rho'}\} \subseteq [n]$ ($1 \leq \rho' \leq \rho$), we introduce a receiver which requests the message x_{i_1} , and has a set $\{x_j : j \in [n] \setminus K\}$ as its side information. It is straightforward to verify that indeed we obtain an instance $(m^*, n, \mathcal{X}^*, f^*)$ satisfying (24). The problem of designing an optimal matrix \mathbf{L} satisfying the (ρ, δ) -Property then becomes equivalent to the problem of finding an optimal (δ, \mathcal{H}^*) -ECIC. Thus, $\mathcal{N}_q[\mathcal{H}^*, \delta]$ is equal to the number of columns in an optimal matrix which satisfies the (ρ, δ) -Property.

The corresponding α -bound and κ -bound for $\mathcal{N}_q[\mathcal{H}^*, \delta]$ can be stated as follows.

Theorem 8.11: Let ρ^* be the smallest number such that a linear $[n, n - \rho^*, \geq \rho + 1]_q$ code exists. Then we have

$$N_q[\rho, 2\delta + 1] \leq \mathcal{N}_q[\mathcal{H}^*, \delta] \leq N_q[\rho^*, 2\delta + 1].$$

Proof: The first inequality follows from the α -bound and from the fact that $\alpha(\mathcal{H}^*) = \rho$, which is due to (24).

For the second inequality, it suffices to show that $\kappa_q(\mathcal{H}^*) = \rho^*$. By Corollary 3.10, an $n \times N$ matrix \mathbf{L} corresponds to an \mathcal{H}^* -IC if and only if $\{\mathbf{L}_i : i \in K\}$ is linearly independent for every $K \in \mathcal{J}(\mathcal{H}^*)$. Since $\mathcal{J}(\mathcal{H}^*)$ is the set of all nonempty subsets of cardinality at most ρ , this is equivalent to saying that every set of at most ρ rows of \mathbf{L} is linearly independent. This condition is equivalent to the condition that \mathbf{L}^T is a parity check matrix of a linear code with the minimum distance at least $\rho + 1$ [26, Chapter 1]. Therefore, a linear \mathcal{H}^* -IC of length N exists if and only if an $[n, n - N, \geq \rho + 1]_q$ linear code exists. Since ρ^* is the smallest number such that an $[n, n - \rho^*, \geq \rho + 1]_q$ code exists, we conclude that $\kappa_q(\mathcal{H}^*) = \rho^*$. ■

Corollary 8.12: The length of an optimal δ -error-correcting linear index code over \mathbb{F}_q which is static under $\Gamma(n, \rho)$ satisfies

$$\mathcal{N}_q[\delta, \mathcal{H}^*] \geq \rho^* + 2\delta,$$

where ρ^* is the smallest number such that an $[n, n - \rho^*, \geq \rho + 1]_q$ code exists.

Proof: This is a straightforward corollary of Theorem 5.1 (the Singleton bound) and Theorem 8.11. ■

Corollary 8.13: For $q \geq \max\{n - 1, \rho + 2\delta - 1\}$, the length of an optimal δ -error-correcting linear index code over \mathbb{F}_q which is static under $\Gamma(n, \rho)$ is $\rho + 2\delta$.

Proof: For $q \geq n - 1$ there exists an $[n, n - \rho^*, \rho + 1]_q$ linear code with $\rho^* = \rho$ (for example, one can take an extended RS code [26, Chapter 11]). Due to Singleton bound, we conclude that $\rho^* = \rho$ is the smallest value such that $[n, n -$

$\rho^*, \rho + 1]_q$ linear code exists. Following the lines of the proof of Theorem 8.11, there exists a δ -error-correcting index code of length $N_q[\rho, 2\delta + 1]$, which is static under $\Gamma(n, \rho)$. As $q \geq \rho + 2\delta - 1$, we have

$$N_q[\rho, 2\delta + 1] = \rho + 2\delta$$

(for example, by taking an extended RS code). Due to Corollary 8.12, this static error-correcting index code is optimal. ■

Remark 8.14: We observe from the proof of Theorem 8.11 that the problem of constructing an optimal linear (non-error-correcting) index code, which is static under $\Gamma(n, \rho)$, is, in fact, equivalent to the problem of constructing a parity check matrix of a classical linear error-correcting code.

Example 8.15: Let $n = 20$, $\rho = 10$, $\delta = 1$ and $q = 2$. From [18], the smallest possible dimension of a binary linear code of length 20 and minimum distance 11 is 3. We obtain that $\rho^* = 17$. We also have $N_2[17, 3] = 22$. Theorem 8.11 implies the existence of a one-error-correcting binary index code of length 22 which can be used for any instance of IC problem, in which each receiver owns at least 10 out of (at most) 20 messages, as side information. It also implies that the length of any such static error-correcting index code is at least $N_2[10, 3] = 14$. Corollary 8.12 provides a better lower bound on the minimum length, which is $17 + 2 = 19$.

Example 8.16: Below we show that with the same number of inputs N and outputs n , a weakly resilient function may have strictly higher resiliency t . From Example 8.15, there exists a linear vectorial Boolean function $\mathbf{f} : (\mathbb{F}_2)^{22} \rightarrow (\mathbb{F}_2)^{20}$ which is 10-weakly 2-resilient. According to [18], an optimal linear $[22, 20]_2$ code has minimum distance $d = 2$. Hence, due to Theorem 8.6, the resiliency of any linear vectorial Boolean function $\mathbf{g} : (\mathbb{F}_2)^{22} \rightarrow (\mathbb{F}_2)^{20}$ cannot exceed one.

The problem of constructing an $n \times N$ matrix \mathbf{L} which satisfies the (ρ, δ) -Property is a natural generalization of the problem of constructing the parity check matrix \mathbf{H} of a linear $[n, k, d \geq \rho + 1]_q$ code. Indeed, \mathbf{H} is a parity check matrix of an $[n, k, d \geq \rho + 1]_q$ code if and only if every set of ρ columns of \mathbf{H} is linearly independent. Equivalently, any nontrivial linear combination of at most ρ columns of \mathbf{H} has weight at least one. For comparison, \mathbf{L} satisfies the (ρ, δ) -Property if and only if any nontrivial linear combination of at most ρ columns of \mathbf{L}^T has weight at least $2\delta + 1$.

Some classical methods for deriving bounds on the parameters of error-correcting codes can be generalized to the case of linear static error-correcting index codes. Below we present a Gilbert-Varshamov-like bound.

Theorem 8.17: Let $V_q(N, 2\delta)$ denotes the volume of q -ary sphere of radius 2δ in \mathbb{F}_q^N given by (9). If

$$\sum_{i=0}^{\rho-1} \binom{n-1}{i} (q-1)^i < \frac{q^N}{V_q(N, 2\delta)},$$

then there exists an $n \times N$ matrix \mathbf{L} which satisfies the (ρ, δ) -Property.

Proof: We build up the set \mathcal{R} of rows of \mathbf{L} one by one. The first row can be any vector in \mathbb{F}_q^N of weight at least $2\delta + 1$. Now suppose we have chosen r rows so that no nontrivial linear combination of at most ρ among these r rows have weight less than $2\delta + 1$. There are at most

$$V_q(N, 2\delta) \sum_{i=0}^{\rho-1} \binom{r}{i} (q-1)^i$$

vectors which are at distance less than $2\delta + 1$ from any linear combination of at most $\rho - 1$ among r chosen rows (this includes vectors at distance less than $2\delta + 1$ from $\mathbf{0}$). If this quantity is smaller than q^N , then we can add another row to the set \mathcal{R} so that no nontrivial linear combination of at most ρ rows in \mathcal{R} has weight less than $2\delta + 1$. The claim follows if we replace r by $n - 1$. ■

Remark 8.18: If we apply Theorem 6.1 to the instance $(m^*, n, \mathcal{X}^*, f^*)$ defined in the beginning of this section, then we obtain a bound, which is somewhat weaker than its counterpart in Theorem 8.17, namely the $n \times N$ matrix \mathbf{L} as above exists if

$$\sum_{i=1}^{\rho} \binom{n}{i} q^{i-1} < \frac{q^N}{V_q(N, 2\delta)}.$$

IX. CONCLUSIONS

In this work, we generalize Index Coding with Side Information problem towards a setup with errors. Under this setup, each receiver should be able to recover its desired message even if a certain amount of errors happen in the transmitted data. This is the first work that considers such a problem.

A number of bounds on the length of an optimal error-correcting index code are constructed. As it is shown in Example 4.8, a separation of error-correcting code and index code sometimes leads to a non-optimal scheme. This raises a question of designing coding schemes in which the two layers are treated as a whole. Therefore, the question of constructing error-correcting index codes with good parameters is still open.

A general decoding procedure for linear error-correcting index codes is discussed. The difference between decoding of a classical error-correcting code and decoding of an error-correcting index code is that in the latter case, each receiver does not require a complete knowledge of the error vector. This difference may help to ease the decoding process. Finding an efficient decoding method for error-correcting index codes (together with their corresponding constructions) is also still an open problem.

The notion of error-correcting index code is further generalized to static index code. The latter is designed to serve a family of instances of error-correcting index coding problem. The problem of designing an optimal static ECIC is studied, and several bounds on the length of such codes are presented.

X. ACKNOWLEDGEMENTS

The authors would like to thank the authors of [7] for providing a preprint of their paper. This work is supported

by the National Research Foundation of Singapore (Research Grant NRF-CRP2-2007-03).

APPENDIX

Lemma A.1: If $\mathcal{G}_{\mathcal{H}}$ is symmetric, then the generalized independence number of \mathcal{H} is the independence number of $\mathcal{G}_{\mathcal{H}}$.

Proof: It suffices to show that if $\mathcal{G}_{\mathcal{H}}$ is symmetric, then the set of generalized independent sets of \mathcal{H} and the set of independent sets of $\mathcal{G}_{\mathcal{H}}$ coincide.

Let H be a generalized independent set in \mathcal{H} . If $|H| = 1$, then obviously H is an independent set in $\mathcal{G}_{\mathcal{H}}$. Assume that $|H| \geq 2$. For any pair of vertices $i, j \in H$, the set $\{i, j\}$ belongs to $\mathcal{J}(\mathcal{H})$. By definition of $\mathcal{J}(\mathcal{H})$, either there is no edge from i to j , or there is no edge from j to i , in $\mathcal{G}_{\mathcal{H}}$. Since $\mathcal{G}_{\mathcal{H}}$ is symmetric, there are no edges between i and j , in neither directions. Therefore, H is an independent set in $\mathcal{G}_{\mathcal{H}}$.

Conversely, let H be an independent set in $\mathcal{G}_{\mathcal{H}}$. For each $i \in H$, since there are no edges from i to all other vertices in H , we deduce that $H \setminus \{i\} \subseteq \mathcal{V}_i$. Due to (3), every subset of H which contains i belongs to $\mathcal{J}(\mathcal{H})$. This holds for an arbitrary $i \in H$. Therefore, every nonempty subset of H belong to $\mathcal{J}(\mathcal{H})$. We obtain that H is a generalized independent set of \mathcal{H} . ■

REFERENCES

- [1] Y. Birk and T. Kol, "Informed-source coding-on-demand (ISCOD) over broadcast channels," in *Proc. IEEE Conf. on Comput. Commun. (INFOCOM)*, San Francisco, CA, 1998, pp. 1257–1264.
- [2] —, "Coding-on-demand by an informed source (ISCOD) for efficient broadcast of different supplemental data to caching clients," *IEEE Trans. Inform. Theory*, vol. 52, no. 6, pp. 2825–2830, 2006.
- [3] S. El Rouayheb, A. Sprintson, and C. Georghiades, "On the index coding problem and its relation to network coding and matroid theory," submitted to *IEEE Trans. Inform. Theory*.
- [4] S. Katti, H. Rahul, W. Hu, D. Katabi, M. Médard, and J. Crowcroft, "Xors in the air: Practical wireless network coding," in *Proc. ACM SIGCOMM*, 2006, pp. 243–254.
- [5] S. Katti, D. Katabi, H. Balakrishnan, and M. Médard, "Symbol-level network coding for wireless mesh networks," *ACM SIGCOMM Comput. Commun. Review*, vol. 38, no. 4, pp. 401–412, 2008.
- [6] Z. Bar-Yossef, Z. Birk, T. S. Jayram, and T. Kol, "Index coding with side information," in *Proc. 47th Annu. IEEE Symp. on Found. of Comput. Sci. (FOCS)*, 2006, pp. 197–206.
- [7] —, "Index coding with side information," *IEEE Trans. Inform. Theory*, to appear.
- [8] E. Lubetzky and U. Stav, "Non-linear index coding outperforming the linear optimum," *Proc. 48th Annu. IEEE Symp. on Found. of Comput. Sci. (FOCS)*, pp. 161–168, 2007.
- [9] Y. Wu, J. Padhye, R. Chandra, V. Padmanabhan, and P. A. Chou, "The local mixing problem," in *Proc. Inform. Theory and Applicat. Workshop*, San Diego, CA, 2006.
- [10] S. El Rouayheb, M. A. R. Chaudhry, and A. Sprintson, "On the minimum number of transmissions in single-hop wireless coding networks," in *Proc. IEEE Inform. Theory Workshop (ITW)*, 2007, pp. 120–125.
- [11] S. El Rouayheb, A. Sprintson, and C. Georghiades, "On the relation between the index coding and the network coding problems," in *Proc. IEEE Symp. on Inform. Theory (ISIT)*, Toronto, Canada, 2008, pp. 1823–1827.
- [12] M. A. R. Chaudhry and A. Sprintson, "Efficient algorithms for index coding," in *Proc. IEEE Conf. on Comput. Commun. (INFOCOM)*, 2008, pp. 1–4.
- [13] N. Alon, A. Hassidim, E. Lubetzky, U. Stav, and A. Weinstein, "Broadcasting with side information," in *Proc. 49th Annu. IEEE Symp. on Found. of Comput. Sci. (FOCS)*, 2008, pp. 823–832.
- [14] R. Ahlswede, N. Cai, S. Y. R. Li, and R. W. Yeung, "Network information flow," *IEEE Trans. Inform. Theory*, vol. 46, pp. 1204–1216, 2000.

- [15] R. Koetter and M. Médard, "An algebraic approach to network coding," *IEEE/ACM Trans. Netw.*, vol. 11, pp. 782–795, 2003.
- [16] W. Haemers, "An upper bound for the shannon capacity of a graph," *Algebr. Methods Graph Theory*, vol. 25, pp. 267–272, 1978.
- [17] S. H. Dau, V. Skachek, and Y. M. Chee, "On the security of index coding with side information," submitted. Also available online at <http://arxiv.org/abs/1102.2797>.
- [18] M. Grassl, "Bounds on the minimum distance of linear codes and quantum codes," available online at <http://www.codetables.de>.
- [19] M. Chudnovsky, N. Robertson, P. Seymour, and R. Thomas, "The strong perfect graph theorem," *Annals of Mathematics*, vol. 164, pp. 51–229, 2006.
- [20] E. R. Berlekamp, R. J. McEliece, and H. C. A. van Tilborg, "On the inherent intractability of certain coding problems," *IEEE Trans. Inform. Theory*, vol. IT-24, no. 3, pp. 384–386, 1978.
- [21] B. Chor, O. Goldreich, J. Håstad, J. Freidmann, S. Rudich, and R. Smolensky, "The bit extraction problem or t-resilient functions," in *Proc. 26th Annu. IEEE Symp. on Found. of Comput. Sci. (FOCS)*, 1985, pp. 396–407.
- [22] C. H. Bennet, G. Brassard, and J. M. Robert, "Privacy amplification by public discussion," *SIAM J. Computing*, vol. 17, pp. 210–229, 1988.
- [23] C. Carlet, *Vectorial Boolean Functions for Cryptography*, ser. Boolean Models and Methods in Mathematics, Computer Science and Engineering. Cambridge University Press, 2010, ch. 9.
- [24] X.-M. Zhang and Y. Zheng, "On nonlinear resilient functions," in *Proc. 14th Annu. Int. Conf. on Theory and Appl. of Cryptographic Tech. (EUROCRYPT)*, 1995, pp. 274–288.
- [25] K. Gupta and P. Sarkar, "Improved construction of nonlinear resilient s-boxes," *IEEE Trans. Inform. Theory*, vol. 51, no. 1, pp. 339–348, 2005.
- [26] F. J. MacWilliams and N. J. A. Sloane, *The Theory of Error-Correcting Codes*. Amsterdam: North-Holland, 1977.